

I claim:

1. A data storage method comprising:
receiving a data set from a client;
5 defining a virtual device to include device portions on a plurality of network servers;
parsing said data set into a plurality of data portions; and
writing each of said data portions to a corresponding one of said device portions.
2. A data storage method according to Claim 1, wherein:
10 said data set is received from said client via a first network; and
and said data portions are written to said device portions via a second network.
3. A data storage method according to Claim 1, further comprising:
receiving additional data sets from additional clients;
15 defining additional virtual devices to include device portions on said plurality of
network servers;
parsing each said additional data set into a plurality of data portions; and
writing each of said data portions of said additional data sets to a corresponding one
of said device portions of said additional virtual devices.
- 20 4. A data storage method according to Claim 3, wherein no more than one user data file
is written to each of said virtual devices.
5. A data storage method according to Claim 3, wherein no more than one directory data
25 file is written to each of said virtual devices.
6. A data storage method according to Claim 3, wherein no more than one meta-data file
is written to each of said virtual devices.

7. A data storage method according to Claim 3, wherein each of said virtual devices contains no more than one type of data.

8. A data storage method according to Claim 1, further comprising incorporating said
5 virtual device into a data structure distributed across said plurality of network servers.

9. A data storage method according to Claim 8, wherein said data set includes new directory data, and incorporating said device into said data structure includes:

parsing a new sub-directory entry;

10 writing said parsed subdirectory entry to a virtual directory device defined across said plurality of network servers, said virtual directory device corresponding to a current directory in which the new directory will be included;

parsing a new meta-data entry;

writing said parsed meta-data entry to a virtual meta-data device defined across said plurality of network servers, said virtual meta-data device corresponding to said current directory;

defining a new virtual directory device across said plurality of network servers;

parsing said new directory data;

writing said parsed new directory data to said new virtual directory device; and

20 defining a new virtual meta-data device across said plurality of network servers.

10. A data storage method according to Claim 8, wherein said data set includes new user file data, and incorporating said virtual device into said data structure comprises:

parsing a new file entry;

writing said parsed file entry to a virtual directory device defined across said plurality of network servers, said virtual directory device corresponding to a current directory in which the new file will be included;

parsing a new meta-data entry;

writing said parsed meta-data entry to a virtual meta-data device defined across said plurality of network servers, said virtual meta-data device corresponding to said current directory.

11. A data storage method according to Claim 1, wherein said step of defining said virtual device comprises:

determining the number of data portions into which said data set is to be parsed;

selecting a number of servers from said plurality of servers corresponding to said number of data portions; and

defining a data portion file for each selected server to store a corresponding one of said data portions.

12. A data storage method according to Claim 11, wherein said number of data portions depends on the type of data in said data set.

13. A data storage method according to Claim 11, wherein each said data portion file is assigned a name, said name including:

an identifier uniquely identifying said virtual device;

a file number uniquely identifying said data portion file with respect to other data portion files corresponding to said virtual device; and

the total number of said data portion files corresponding to said virtual device.

14. A data storage method according to Claim 11, wherein:

said step of defining said virtual device includes defining one of said device portions to include parity data; and

said step of parsing said data includes generating said parity data.

5

15. A data storage method according to Claim 11, wherein said step of selecting a number of servers from said plurality of servers includes selecting said servers based at least in part on the available storage capacity of said servers.

10

16. A data storage method according to Claim 11, wherein said step of writing each of said data portions to a corresponding one of said device portions includes transmitting each of said data portions to a corresponding one of said network servers as a corresponding one of said data portion files.

1003604-1200
T0022T "S403E00T

17. A data storage method according to Claim 16, further comprising transmitting a signal to said client indicating that said data set has been stored, after said data set has been received, but before said data set has been written to said virtual device.

20

18. A data storage method according to Claim 16, further comprising:

writing said data set to local non-volatile memory; and

transmitting a signal to said client indicating that said data set has been stored, after said data set has been written to said non-volatile memory, but before said data set has been written to said virtual device.

25

19. A data storage method according to Claim 16, further comprising transmitting a signal to said client indicating that said data has been stored, only after said data set has been written to said virtual device.

20. A data storage method according to Claim 16, further comprising transmitting a signal to said client indicating that said data has been stored, said signal being transmitted at one of the following times depending on a predetermined criteria:

after said data set has been received, but before said data set has been written to said virtual device;

after said data set has been written to local non-volatile memory, but before said data set has been written to said virtual device; or

only after said data set has been written to said virtual device.

21. A data storage method according to Claim 20, wherein said predetermined criteria includes a data type of said data set.

22. A data storage method according to Claim 20, wherein said predetermined criteria includes a file name extension associated with said data set.

23. A data storage method according to Claim 16, further comprising transmitting a commit signal to each of said corresponding network servers, said commit signals causing said network servers to commit said data portions to memory.

24. A data storage method according to Claim 23, further comprising:
determining whether a confirmation signal has been received from each of said corresponding network servers, said confirmation signals indicating that said network servers have committed said data portions to memory; and
writing a write failure entry to at least one fact server, said write failure entry identifying any of said corresponding network servers from which said confirmation signals are not received.

25. A data storage method according to Claim 24, wherein said write failure entry is written to at least two fact servers.

26. A data storage method according to Claim 25, wherein said fact servers each reside on a different one of said plurality of network servers.

27. A data storage method according to Claim 25, further comprising:

periodically polling said fact servers;
correcting any data portion files corresponding to incomplete writes identified by said fact servers.

28. A data storage method according to Claim 23, further comprising:

transmitting a ready signal to a backup controller after transmitting each of said data portions to a corresponding one of said network servers; and
transmitting a done signal to said backup controller after transmitting said commit signals to said servers;
whereby said backup controller, responsive to receiving said ready signal and not receiving said done signal, transmits a commit signal to said corresponding network servers.

29. A data storage method according to Claim 28, further comprising:

determining whether a confirmation signal has been received from each of said corresponding network servers, said confirmation signals indicating that said network servers have committed said data portions to memory; and
writing a write failure entry to at least one fact server, said write failure entry identifying any of said corresponding network servers from which said confirmation signals are not received; and wherein
said backup controller, responsive to receiving said ready signal and not receiving said done signal, performs said steps of determining whether said confirmation signals have been received and writing said write failure entry to said fact server.

30. A data storage method according to Claim 29, wherein said write failure entry is

written to at least two fact servers.

31. A data storage method according to Claim 30, wherein said fact servers each reside on a different one of said plurality of network servers.

5 32. A data storage method according to Claim 1, wherein said step of receiving said data set from said client comprises:

writing said data set to local nonvolatile data storage; and

writing a local data entry to a fact server, said local data entry indicating that valid data is stored in said local nonvolatile data storage.

10 33. A data storage method according to Claim 29, wherein said local data entry is written to at least two fact servers.

15 34. A data storage method according to Claim 33, wherein said fact servers each reside on a different one of said plurality of network servers.

20 35. A data storage method according to Claim 32, further comprising:

removing said data set from said local nonvolatile memory after said data set is

written to said virtual device; and

updating said fact server to indicate that said data set is no longer in said local nonvolatile data storage.

25 36. A data retrieval method comprising:

receiving a data request from a client;

retrieving a virtual device definition corresponding to the requested data, said virtual device definition identifying device portions located on a plurality of network servers;

retrieving data portion files from said device portions;

collating said retrieved data portion files to generate the requested data; and

30 transmitting the requested data to said client.

37. A data retrieval method according to Claim 36, wherein said step of retrieving said virtual device definition comprises:

retrieving virtual meta-data device information from a current directory;
retrieving meta-data portion files from said plurality of network servers;
collating said meta-data portion files to generate meta-data; and
retrieving said virtual device definition from said meta-data.

38. A data retrieval method according to Claim 36, wherein said step of retrieving data portion files from said device portions comprises:

transmitting requests for said data portion files to network servers corresponding to said device portions; and
receiving said data portion files from said network servers.

39. A data retrieval method according to Claim 38, wherein each said data portion file is assigned a name, said name including:

an identifier uniquely identifying the requested data;
a file number uniquely identifying said data portion file with respect to other data portion files corresponding to the data; and
the total number of said data portion files corresponding to the requested data.

40. A data retrieval method according to Claim 38, wherein said step of retrieving data portion files from said device portions comprises:

determining which one of a plurality of controllers has access to said virtual device, said controllers residing on said network servers; and
invoking said controller with access to said virtual device to retrieve said data portion files.

41. A data retrieval method according to Claim 40, wherein said step of retrieving said virtual device definition comprises:

5 determining which one of said plurality of controllers has access to another virtual device storing said virtual device definition; and
 invoking said controller with access to said another virtual device to retrieve said virtual device definition.

10 42. A data retrieval method according to Claim 41, wherein said steps of determining which ones of said plurality of controllers have access to said virtual devices include selecting said controllers depending on the type of data stored in said virtual devices.

15 43. A data retrieval method according to Claim 38, wherein:
 said step of receiving said data portion files from said network servers includes receiving all but one of said data portion files; and
 said step of collating said data portion files includes generating said one data portion file based on parity data included in said received data portion files.

20 44. A data retrieval method according to Claim 36, further comprising:
 polling at least one fact server to determine whether said virtual device includes a potentially corrupt data portion file; and
 reconstructing said potentially corrupt data portion file.

45. A data retrieval method according to Claim 36, further comprising:

polling at least one fact server to determine whether said virtual device includes a
potentially corrupt data portion file;

polling at least one fact server to determine whether said requested data is stored in
nonvolatile data storage; and

retrieving said requested data from said nonvolatile data storage instead of said virtual
device, if said virtual device includes a potentially corrupt data portion file and
said requested data is stored in said nonvolatile data storage.

46. A data retrieval method according to Claim 36, wherein:

communication with said client is via a first network; and

communication with said network servers is via a second network.

47. A data retrieval method according to Claim 46, wherein communication between said
client and said network servers is via said first network.

48. A computer-readable medium having code embodied therein for causing an
electronic device to perform the method of Claim 1.

49. A computer-readable medium having code embodied therein for causing an
electronic device to perform the method of Claim 2.

50. A computer-readable medium having code embodied therein for causing an
electronic device to perform the method of Claim 3.

51. A computer-readable medium having code embodied therein for causing an
electronic device to perform the method of Claim 4.

52. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 5.

5 53. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 6.

54. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 7.

10 55. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 8.

56. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 9.

57. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 10.

20 58. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 11.

59. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 12.

25 60. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 13.

61. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 14.

30

62. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 15.

5 63. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 16.

64. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 17.

10 65. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 18.

66. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 19.

67. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 20.

20 68. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 21.

69. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 22.

25 70. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 23.

71. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 24.

30

72. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 25.

73. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 26.

74. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 27.

75. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 28.

76. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 29.

77. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 30.

78. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 31.

79. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 32.

80. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 33.

81. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 34.

82. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 35.

5 83. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 36.

84. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 37.

10 85. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 38.

86. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 39.

87. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 40.

20 88. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 41.

89. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 42.

25 90. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 43.

91. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 44.

30

92. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 45.

93. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 46.

94. A computer-readable medium having code embodied therein for causing an electronic device to perform the method of Claim 47.

95. A data storage system comprising:

- a network interface to facilitate communication with clients and with a plurality of network servers; and
- a file server, responsive to receiving a file to be stored from one of said clients, and operative to define a virtual device to include device portions on said plurality of network servers, to parse said file into a plurality of file portions, and to write each of said file portions to a corresponding one of said device portions.

96. A data storage system according to Claim 95, wherein said network interface comprises:

- a first network adapter to facilitate communication between said data storage system and said clients; and
- a second network adapter to facilitate communication between said data storage system and said plurality of network servers.

97. A data storage system according to Claim 95, further comprising a data storage device for storing at least one of said file portions, whereby said data storage system is capable of functioning as one of said plurality of network servers.

98. A data storage system according to Claim 95, wherein said virtual device is limited to storing no more than one file.

99. A data storage system according to Claim 95, wherein said file server includes:
a client process for receiving said file from said client; and
a distribution controller, responsive to said client process, and operative to determine
the number of file portions into which said file is to be parsed, to select a number
of servers from said plurality of servers corresponding to said number of file
portions, and to define a portion file for each selected server to store a
corresponding one of said file portions.

100. A data storage system according to Claim 99, wherein said distribution controller
determines the number of file portions into which the file is to be parsed based, at least in part,
on the file type of the file.

101. A data storage system according to Claim 99, wherein said distribution controller is
further operative to assign a name to each of said portion files, said name including:
an identifier uniquely identifying said file;
a file number uniquely identifying said portion file with respect to other portion files
associated with said data set; and
the total number of said portion files associated with said data set.

102. A data storage system according to Claim 99, wherein said distribution controller is
further operative to:
define an additional portion file to include parity data; and
to generate said additional portion file from said portion files.

103. A data storage system according to Claim 99, wherein said distribution controller is further operative to:

- to determine the available storage on each of said plurality of servers; and
- to select servers from said plurality of servers based at least in part on the available storage capacity of said servers.

104. A data storage system according to Claim 99, further comprising an input/output process, responsive to receiving requests for locally stored portion files from said distribution controller, and operative to retrieve and transmit said portion files to said distribution controller.

105. A data storage system according to Claim 104, wherein said distribution controller is further operative to transmit a signal to said client indicating that said file has been stored, said signal being transmitted after said file has been received by said client process, but before said file has been written to said virtual device.

106. A data storage system according to Claim 104, wherein said distribution controller is further operative to:

- write said file to local non-volatile memory; and
- to transmit a signal to said client indicating that said file has been stored, after said file has been written to said non-volatile memory, but before said file has been written to said virtual device.

107. A data storage system according to Claim 104, wherein said distribution controller is further operative to transmit a signal to said client indicating that said file has been stored, only after said file has been written to said virtual device.

108. A data storage system according to Claim 104, wherein said distribution controller, responsive to predetermined criteria, is further operative to transmit a signal to said client indicating that said file has been stored, said signal being transmitted at one of the following
5 times:

after said file has been received, but before said file has been written to said virtual device;

after said file has been written to local non-volatile memory, but before said file has been written to said virtual device; or

only after receiving confirmation from said distribution controller that said file has been written to said virtual device.

109. A data storage system according to Claim 108, wherein said predetermined criteria includes a file type of said file.

110. A data storage system according to Claim 109, wherein said predetermined criteria includes a file name extension associated with said file.

111. A data storage system according to Claim 104, wherein said distribution controller is further operative to:

determine whether a confirmation signal has been received from each of said corresponding network servers, said confirmation signals indicating that said network servers have committed said portion files to memory; and
write a write failure entry to at least one fact server, said write failure entry identifying
25 any of said corresponding network servers from which said confirmation signals are not received.

112. A data storage system according to Claim 111, wherein said distribution controller is operative to write said write failure entry to at least two fact servers.

113. A data storage system according to Claim 112, wherein:

said fact servers each reside on a different one of said plurality of network servers;

and

said data storage system further includes a fact server to receive entries from said
network servers.

114. A data storage system according to Claim 112, further comprising:

a local controller operative to poll said fact servers; and wherein,

responsive to a signal from said local controller, said distribution controller is

operative to correct any portion files corresponding to incomplete writes identified
by said fact servers.

115. A data storage system according to Claim 104, wherein said distribution controller
is further operative to transmit a commit signal to each of said corresponding network servers,
said commit signals causing said network servers to commit said portion files to memory.

116. A data storage system according to Claim 115, wherein said distribution controller
is further operative to:

transmit a ready signal to a backup controller after transmitting each of said portion
files to a corresponding one of said network servers; and

transmit a done signal to said backup controller after transmitting said commit signals
to said servers;

whereby said backup controller, responsive to receiving said ready signal and not
receiving said done signal, transmits a commit signal to said corresponding
network servers.

117. A data storage system according to Claim 116, wherein said distribution controller is further operative to:

receive confirmation signals from each of said corresponding network servers, said confirmation signals indicating that said network servers have committed said portion files to memory; and

write a write failure entry to at least one fact server, said write failure entry identifying any of said corresponding network servers from which said confirmation signals are not received.

118. A data storage system according to Claim 117, wherein said distribution controller is operative to write said write failure entries to at least two fact servers.

119. A data storage system according to Claim 118, wherein said fact servers each reside on a different one of said plurality of network servers.

120. A data storage system according to Claim 119, further comprising:

a local controller operative to poll said fact servers; and wherein

responsive to a signal from said local controller, said distribution controller is

operative to correct any portion files corresponding to a write failure entry in at least one of said fact servers.

121. A data storage system according to Claim 95, wherein said file server is operative to:

write said file to local nonvolatile data storage pending said write of said file portions to said device portions; and

writing a local data entry to a fact server, said local data entry indicating that said file is stored in said local nonvolatile data storage.

122. A data storage system according to Claim 121, wherein said local data entry is written to at least two fact servers.

123. A data storage system according to Claim 122, wherein said fact servers each reside
5 on a different one of said plurality of network servers.

124. A data storage system according to Claim 121, wherein said file server is further operative to:

remove said file from said local nonvolatile memory after said file portions are
10 written to said device portions; and
update said fact server to indicate that said file is no longer in said local nonvolatile data storage.

125. A data storage system comprising:

a client interface operative to receive a file request from a client;
a file server responsive to said file request, and operative to retrieve a virtual device
definition corresponding to said requested file, said virtual device definition
identifying file portions located on a plurality of network servers, to retrieve said
file portions, and to collate said retrieved file portions to generate said requested
20 file.

126. A data storage system according to Claim 125, wherein said client interface is further operative to transmit said requested file to said client.

127. A data storage system according to Claim 126, wherein said file server comprises:
a local controller responsive to said file request, and operative to retrieve virtual meta-data device information associated with said requested file from a current directory, said meta-data device information identifying meta-data device portions on said plurality of network servers;
a distribution controller responsive to said virtual meta-data device information, and operative to retrieve meta-data portion files from said meta-data device portions, to collate said meta-data portion files to generate meta data, said meta-data including said virtual device definition corresponding to said requested file.

128. A data storage system according to Claim 126, wherein said file server comprises a distribution controller responsive to said virtual device definition and operative to transmit requests for said file portions to said servers, and to receive said file portions from said servers.

129. A data storage system according to Claim 128, wherein said file portions are identified by a file portion name, said file portion name including:
an identifier uniquely identifying said requested file;
a second identifier uniquely identifying said file portion with respect to said other file portions corresponding to said requested file; and
the total number of file portions corresponding to said requested file.

130. A data storage system according to Claim 125, wherein said client interface, responsive to said file request, is operative to determine which of a plurality of controllers has access to said virtual device.

131. A data storage system according to Claim 130, further comprising:
at least one of said plurality of controllers; and wherein
at least one other of said plurality of controllers is located on one of said network servers of said plurality of network servers.

132. A data storage system according to Claim 130, wherein which controller has access to said virtual device depends on the type of file stored in said virtual device.

133. A data storage system according to Claim 132, wherein each of said controllers is
5 able to handle only one type of file.

134. A data storage system according to Claim 125, wherein said file server, responsive to receiving all but one of said file portions, is operative to generate said one file portion from parity data in one of said received file portions.

10

135. A data storage system according to Claim 125, wherein:
said client interface communicates with clients via a first network; and
said file server communicates with said plurality of network servers via a second
network.